高頻度データを用いたニュースと韓国株式の分析:ニュース分類の条件探索 Analysis of News and Stock Returns in Korean Markets Using High Frequency Data: Exploring Conditions of Classification Algorithm

> 尹聖在^{*1} Sungjae Yoon

菅愛子*1 Aiko Suge 高橋大志^{*1} Hiroshi Takahashi

*1 慶應義塾大学大学院経営管理研究科 Graduate School of Business Administration, Keio University

News articles play an important role in financial markets. This study analyzes the relationship between news articles and stock price fluctuations using high frequency trading data in Korean stock markets. Especially, we analyze differences in market reactions according to languages of news articles. In order to understand the influences of news articles, this study explores conditions of Long Short Term Memory (LSTM) models that classify news articles.

1. はじめに

フィンテックなど金融分野における技術の進展を背景とし, 資 産運用への関心が高まっている. 資産運用分野において, 新た な情報技術を用いた分析などが盛んに報告される中, 金融市 場において発信されるニュースは, 資産価格に関する重要な情 報の一つに挙げられる. 効率的な市場において, 情報は資産 価格に速やかに反映されるという理論背景において, 日々発信 されるニュースと株式市場の関連性について, これまで数多くの 実証研究がなされてきた. 近年では, ニュース記事の内容を定 量化・指標化する試みがなされ, それらと投資家行動に及ぼす 要素とを関連付ける研究もされている.

我々は以前,様々な言語のニュースが株価に与える影響を 韓国株式市場において分析した(図 1)[Yoon, 2017] [Yoon, 2018]. 更にその中で,Long Short Term Memory (LSTM)ブロッ ク [Hochreiter & Schmidhuber, 1995] [Hochreiter & Schmidhuber, 1997]を用いた深層学習モデルにて,英語ニュー スによる韓国株式のリターン分類を報告した.そのモデル作成 時のハイパーパラメータ探索の試みを本稿で報告する.



2. 目的

機械学習モデルの学習において、良いモデルを作るために ハイパーパラメータ(外部から与えるモデルのパラメータ)のチュ ーニングが必要となる。良いモデルの一例として、精度の高い モデルが挙げられる。目的の精度をもとめてモデル学習を繰返 す際、選択するアルゴリズム、データ量、計算環境などによって、

連絡先:菅愛子,慶應義塾大学大学院経営管理研究科, aikosuge@keio.jp



図2:モデルの概要

一回の試行に時間を要することも多い. そこで効率よく探索を行うため,様々な試みがなされており,ランダムサーチ,グリッドサーチ,ベイズ最適化,などが知られている.

深層学習では、多数のパラメータが存在する. 例えば筆者ら 用いたモデルは図 2 のように要約されるが、そこで用いるハイパ ーパラメータとして例えば、学習係数、隠れ層のユニット数, 層 の段数, epoch 数, dropout, batch size などが挙げられる.

本稿ではパラメータ探索の始めとして、どの段階におけるパラ メータの影響が大きいのか、ベクトル(入力段階)、分類学習の 段階、で大きく分けパラメータの変更を行い、分類結果(精度) の変化を観察した.

3. データ

3.1 ニュースデータ

ニュースデータは韓国株式市場に上場している時価総額上 位5社(2012年末時点)について Thomson Reuters から抽出し た. 2013年ニュースをモデル作成用(学習用)に用い, 2014年 ニュースを予測に用いた.ニュース発信言語は英語の他,韓国 語,日本語,中国語など多岐に渡るが,本稿では英語ニュース を扱う.

3.2 株価データ

Thomson Reuters の高頻度データから、上記で抽出したニュース発信前後の株価を 15 分ごとに取得した.ニュース発信 30 分前に比較し、ニュース発信 30 分後に株式が 0.3%以上値上 がりしたものをポジティブ、0.3%以上値下がりしたものをネガティブ、それら以外をニュートラルなニュースであると分類した.結果、表1のようなデータセットが得られた.

	ネガティブ	ニュートラル	ポジティブ	計			
2013年	15.9%	63.5%	20.6%	908			
2014年	19.3%	62.8%	17.9%	643			

表 1:ニュースデータの分布

4. 方法

基本モデルを設定し、その中のハイパーパラメータを変化さ せ、精度の変化を検証した.

4.1 基本モデルの設定

図 2 に示すようにニュース分類をセットした. すなわち,文章 (単語)のベクトル化には skip-gram 法 [Mikolov, 13] を,株価 文章分類の学習には LSTM ブロックを用いた. 2013 年ニュー スデータ(908 サンプル)を訓練データとテストデータに分け (9:1),モデル学習には訓練データを用いた.

4.2 パラメータの変更

入力段階においては、文章をベクトル化する際に用いている skip-gram のパラメータ、すなわち学習に用いる単語数やベクト ルの次元、学習効率(initial learning rate: initial LR, epoch 数) を変化させた.また、ニュース分類段階では LSTM ブロックの学 習効率(initial LR, epoch 数)を変化させた.

4.3 結果の検証

訓練データとテストデータの,精度の比較を行った.また,い くつかのモデルにおいて,2014年ニュースデータを用い,株価 リターン分析を行い,分類予測を行った.なお,結果は20回の 学習を行い,その精度を平均して示している.

5. 結果

表 2 にパラメータ探索の結果を一部示す. なお, 結果は 20 回の学習を行い, その精度を平均して示している.

	入力側	LSTM ブロック		検証	
条件	単語数: n	Initial LR	Epoch 数	訓練精度	検証精度
1	100	0.01	100	63.6%	63.3%
2	100	0.01	200	63.4%	64.4%
3	100	0.01	300	63.4%	64.4%
4	100	0.1	200	70.7%	63.3%
5	200	0.01	200	63.3%	63.6%

表 2:条件による精度の変動

今回の結果では、入力する単語数 n が精度に与える影響は 小さいことが分かった.しかしながら、今回用いたテキストデータ の最長は 900 単語を超えるものもあり、今後、入力条件につい て詳細に分析していきたい.また、分類学習における initial LR を 0.1 以上にすると、訓練データが示すモデル精度と検証デー タが示す精度で乖離が見られ、過学習が起きていることが分か った.精度が 60%程度にとどまっている原因として、サンプル数 の少なさや、データの不均衡分布が考えられる.今後はその点 を調整して学習モデルの改善に取り組む予定である.

最後に表 2 の条件 2 について, 2014 年のニュースデータを 用い分類を行い,実際の株価リターン(平均)を算出した.その 分類結果を表 3 に示す.この時の予測精度は 55.7%となった. 表 3: 学習モデルによる 2014 年ニュース分類(条件 2)

予測	ネガティブ	ニュートラル	ポジティブ	計
分布	17.0%	73.7%	9.3%	643

今回はデータが限られており,別の期間などのデータを用いて,モデルや精度の検証をする必要がある.

6. まとめ・課題

ニュースによる株価分類を深層学習にてモデル化した際の ハイパーパラメータについて探索した.数多くのパラメータが存 在するが,今回は入力単語数と,モデル学習における initial RL と epoch 数について変動させ,モデル精度に与える影響を 観察した.過学習が起きない一定の条件を見出したものの,全 体の精度は 60%程度にとどまった.今後,更なる精度向上のた めにモデルやパラメータ設定が課題である.

また深層学習以外の分類法との比較も望まれる.

参考文献

- [Hochreiter & Schmidhuber, 95] Sepp Hochreiter and J^{*}urgen Schmidhuber: Long Short Term Memory, Technical Report FKI-207-95, Technische Universit^{*}at M^{*}unchen, M^{*}unchen, August 1995.
- [Hochreiter & Schmidhuber, 97] Sepp Hochreiter and J^{*}urgen Schmidhuber: Long Short Term Memory, Neural Computation, 9(8):1735–1780, November 1997.
- [Mikolov, 13] Mikolov, T., Chen, K., Corrado, G., and Dean, J.: Efficient estimation of word representations in vector space, ICLR, 2013.
- [Yoon, 17] Yoon, S.J., Suge, A., and Takahashi, H.: Analyzing the Influence of News Articles on the Korean Stock Market with High Frequency Trading Data, JSAI International Symposia on AI, Workshop 3: Artificial Intelligence of and for Business (AI-Biz 2017), 2017.
- [Yoon, 18] Yoon, S.J., Suge, A., and Takahashi, H.: The Analysis on Korean Stock Market by High Frequency Data and Text Classification Using LSTM, JSAI Special Interest Group on Business Informatics (SIG-BI #8), 2018.